

## EXHIBIT 1

~~FILED UNDER SEAL~~

PUBLIC REDACTED  
VERSION



	A	B	C	D	E	F	G	H
1	Unique, De-duplicated Hit Counts for all Terms Combined		The Times's Proposed Search Terms	Total Documents Hit by Term Including Families	OpenAI's Initial Search Term Proposal	Total Documents Hit by Term Including Families	131,601	Court's Ruling
2				546,902				
38	27		data w/10 (clean* OR *process* OR curat* OR select* OR augment* OR samp* OR gather* OR copy* OR publish*) (CC NOT w/25 (.com OR (To AND BCC))) w/50 (data OR content OR *train* OR dataset* OR *gpt* OR dv3 OR dv-3 OR davinci-3 OR generative OR SLM OR Bing OR Copilot OR Azure OR genai OR Prometheus OR Orchestrat*) (train* OR stage* OR layer*) w/5 (learn* OR iteration OR cycle OR batch OR epoch) C4 webtext* OR WT OR WT2	71,644	N/A			
39	28			123,352	N/A			
40	29			30,177	N/A			
41	30			13,205	N/A			
42	31			5,792	webtext		7,087	
43	32				commoncrawl OR "common crawl"		7,669	
44	33		(*gpt* OR dv3 OR dv-3 OR davinci-3 OR generative OR SLM OR Bing OR Copilot OR Azure OR genai OR Prometheus OR Orchestrat*) w/50 (journal* OR news* OR report* OR info* OR knowledge OR article* OR story OR stories OR publish* OR headline* OR cover* OR syndicat* OR podcast* OR lede OR scoop* OR editorial OR opinion) (impact* w/5 assess*) w/25 (*gpt* OR dv3 OR dv-3 OR davinci-3 OR genai OR generative OR SLM OR Bing OR Copilot OR Azure)	78,357	(ChatGPT OR GPT-3.5 OR GPT-3.5T OR "GPT-3.5 turbo" OR GPT-4 OR GPT-4T OR "GPT-4 turbo") w/25 (journal* OR news*)		4,255	
45	34			44	N/A			
46	35			2,512	N/A			
47	Section E: Terms Relating to Technical Documentation and Processes							
48	36		(reward OR policy) w/10 model*	18,701	N/A			
49	37		(few OR one OR zero) w/10 shot*	27,764	N/A			
50			(training OR stag* OR layered) w/5 (learn* OR iteration OR cycle OR batch OR epoch OR curve)					
51	38			21,245	N/A			
52	39		"model completion"	705	N/A			
53	40		data w/25 scrap*	5,539	N/A			
54	41		(*validation data* OR "The Pile" OR preprocessing OR "test set" OR "training set" OR "evaluation set" OR "text extraction" OR "semantic analysis" OR "source attribution") w/25 (discriminator OR curricul* OR classif* OR *score* OR frequen* OR contamination* OR duplicat* OR regression* OR cluster* OR feature* OR label* OR attribute* OR categ* OR value* OR variable* OR inference OR distillation OR recipe*)	3,756	N/A			
55	42		(*validation data* OR "The Pile" OR preprocessing OR "test set" OR "training set" OR "evaluation set" OR "text extraction" OR "semantic analysis" OR "source attribution") w/25 (rule* OR guidance* OR weight* OR perplexity OR ROUGE OR F1 OR classif* OR benchmark OR whitelist OR blacklist OR safeguard* OR guardrail* OR quality OR recitation* OR scrap* OR checkpoint)	1,676	N/A			
56	43		KPI OR PRD OR (performance w/25 indicator)	73	N/A			
57	44		KPI OR PRD OR (performance w/25 indicator)	73	N/A			
58	45		measur* w/5 pipe*	130	N/A			
59	46		(model w/25 card) OR "product brief" OR "product requirement" OR "data card" OR datasheet OR "data doc**"	430	N/A			
60	47		(product* OR design* OR technical OR architecture OR system* OR launch*) w/10 (specification* OR requirement* OR brief* OR report* OR chart OR log) embed* w/15 (function OR space OR vector OR index* OR indices)	42,250	N/A			
61	48		relevance OR "query tuning" OR (source* w/5 (trust* OR credib*)) OR (standard* w/5 editor*) OR (criteri* w/5 (exclusion OR inclusion))	9,855	N/A			
62	49			15,904	N/A			
63	50			1,299	N/A			
64	51		Azure AND compute* (journal* OR news* OR report* OR info* OR knowledge OR article* OR story OR stories OR publish* OR headline* OR cover* OR syndicat*) AND ((data* w/5 quality) OR (high w/5 quality))	20,934	Azure AND ("super computer" OR supercomputer)		1,335	
65	52			34,433	dataset* w/10 "high quality"		14,570	
66	Section F: Terms Relating to Collaboration with Microsoft							
67	53			13,675	N/A			
68	54	crescendo		834	N/A			
69	55			167	N/A			
70					(Microsoft OR MSFT) AND ("super computer" OR supercomputer)			
71	56		mail[@microsoft.com] AND (data* OR train* OR *gpt* OR dv3 OR dv-3 OR davinci-3 OR generative OR Genai OR SLM OR Bing OR Copilot OR Azure OR "compute")				2,675	
72			MS w/25 (data* OR train*)					
73	57			13,645	mail[@microsoft.com] AND (train* w/10 (data* OR material* OR source*))		151	
74	58			3,719	N/A			
75	Section G: Search Terms Related to the Market for and Value of Content							
76	59		(*Licens* OR agreement* OR compensat* OR pay OR paid) AND (*gpt* OR dv3 OR dv-3 OR davinci-3 OR genai OR generative OR SLM OR Bing OR Copilot OR Azure OR Prometheus OR Orchestrat*)	46,947	N/A			
77	60		(license OR agreement) AND (train* w/15 model*)	6,127	N/A			
78	61		(public* OR "open source" OR "creative commons" OR government OR 1920 OR 1921 OR 1922 OR 1923 OR 1924 OR 1925 OR 1926 OR 1927 OR 1928) w/25 (*gpt* OR dv3 OR genai OR generative OR SLM OR Bing OR Copilot OR Azure)	16,071	N/A			

	A	B	C	D	E	F	G	H
1	Unique, De-duplicated Hit Counts for all Terms Combined		The Times's Proposed Search Terms	Total Documents Hit by Term Including Families	OpenAI's Initial Search Term Proposal	Total Documents Hit by Term Including Families	131,601	Court's Ruling
2			(MAI OR MAI-1 OR Phi OR Phi-2 OR Phi-3) w/50 (azure OR CPU OR GPU OR compute* OR invest* OR financ* OR profit* OR revenue* OR stake OR money OR process* OR server* OR hardware)	546,902				
79			62		515	N/A (MAI OR MAI-1 OR Phi OR Phi-2 OR Phi-3) AND azure		
80			63				699	
81								
82	Section H: Search Terms OpenAI Seemingly Re-Used from Other Cases							
83								
84			64			Libgen* OR "Library Genesis" "books corpus" OR (book* w/5 corp*) OR bookscorp* Gutenberg z-library OR zlibrary OR zlib "Google books" B-ok "Open Library"	4,748	
85			65				2,721	
86			66				3,111	
87			67				1,287	
88			68				1,698	
89			69				1,077	
90			70				1,520	
91			71			book* w/5 ("data quality" OR "high quality")	174	
92			72			smashwords	1,188	
93			73			Bibliotik	18	
94			74			"Anna's Archive"	11	
95						Sci-Hub OR scihub OR "sci hub"	342	